

Claims

We claim:

1. A video conferencing system comprising:
 - an image pickup device for generating image signals representative of an image;
 - an audio pickup device for generating audio signals representative of sound from an audio source; and
 - a multimodal integration architecture system for processing said image signals and said audio signals to determine a direction of the audio source relative to a reference point.
2. The video conferencing system of claim 1 wherein said multimodal integration architecture system further comprises:
 - an audio source localization system;
 - a computer vision person detection system; and
 - a multimodal speaker detection system.
3. The video conferencing system of claim 2, further comprising an integrated housing for an integrated video conferencing system incorporating the image pickup device, the audio pickup device, and the multimodal integration architecture system.
4. The video conferencing system of claim 3, wherein the integrated housing is sized for

2 being portable.

1 5. The video conferencing system of claim 2, further comprising an electronic pan tilt zoom
2 system for electronically manipulating the image signals to effectively provide at least one of
3 variable pan, tilt, and zoom functions.

1 6. The video conferencing system of claim 5, wherein the image pickup device is a stationary
2 camera.

1 7. The video conferencing system of claim 5, wherein the multimodal integrated architecture
2 system provides control signals to the electronic pan tilt zoom system.

1 8. The video conferencing system of claim 7, wherein the audio source moves relative to
2 the reference point, the audio source localization system detects the movement of the audio
3 source, and, in response to the movement, the audio source localization system causes a change in
4 the field of view of the image pickup device.

1 9. The video conferencing system of claim 5, wherein the audio pickup device is comprised of
2 an array of two microphones.

1 10. A method comprising the steps of:

2 generating, at an image pickup device, image signals representative of an image;

3 generating, at an audio pickup device, audio signals representative of sound from an audio
4 source;

5 processing the image signals and the audio signals to determine a direction of the audio
6 source relative to a reference point;

7 manipulating the image signals to produce refined image signals; and

8 outputting said refined image signals.

11. The method of claim 10 further comprising the steps of:

9 applying said audio signals to an audio source localization system;

10 applying said image signals to a computer vision person detection system;

11 processing said audio signals and said image signals with a multimodal speaker detection
12 system;

13 generating control signals based on the determined direction of the audio source;

14 applying the control signals to an electronic pan tilt zoom system to mimic the effect of at
15 least one function of a movable camera, said function selected from the group consisting panning,
16 tilting, and zooming said movable camera; and

17 providing an output from said electronic pan tilt zoom system.

1 12. The method of claim 10, further comprising electronically varying a field of view of the
2 image pickup device in response to the control signals.

3 13. The method of claim 10, wherein processing the audio signals includes determining an
4 audio based direction of the audio source based on the audio signals.

1 14. The method of claim 12, wherein the audio source moves relative to a reference point,
2 and wherein processing the audio signals further includes:
3 detecting the movement of the audio source; and
4 causing electronically, in response to the movement, an increase in the field of view of the
5 image pickup device.

1 15. The method of claim 12, further comprising the step of supplying control signals, based on
2 the audio based direction, for electronically panning, tilting, or zooming said image pickup
3 device.

1 16. A video conferencing system comprising:
2 two microphones for generating audio signals representative of sound from a speaker;
3 a video camera for generating video signals representative of a video image;
4 an electronic pan tilt zoom system for manipulating video images to produce the visual
5 effects of panning, tilting, and/or zooming;
6 a processor for processing the video signals and the audio signals to determine a direction
7 of a speaker relative to a reference point and supplying control signals to the electronic pan tilt
8 zoom system for producing images that include the speaker in the field of view of the camera, the
9 control signals being generated based on the determined direction of the speaker; and
a transmitter for transmitting audio and video signals for video conferencing.